

PYTHON REGULAR EXPRESSIONS

A *regular expression* is a special sequence of characters that helps you match or find other strings or sets of strings, using a specialized syntax held in a pattern. Regular expressions are widely used in UNIX world.

The module **re** provides full support for Perl-like regular expressions in Python. The re module raises the exception `re.error` if an error occurs while compiling or using a regular expression.

We would cover two important functions, which would be used to handle regular expressions. But a small thing first: There are various characters, which would have special meaning when they are used in regular expression. To avoid any confusion while dealing with regular expressions, we would use Raw Strings as **r'expression'**.

The *match* Function

This function attempts to match RE *pattern* to *string* with optional *flags*.

Here is the syntax for this function:

```
re.match(pattern, string, flags=0)
```

Here is the description of the parameters:

Parameter	Description
pattern	This is the regular expression to be matched.
string	This is the string, which would be searched to match the pattern at the beginning of string.
flags	You can specify different flags using bitwise OR (). These are modifiers, which are listed in the table below.

The `re.match` function returns a **match** object on success, **None** on failure. We would use `group(num)` or `groups()` function of **match** object to get matched expression.

Match Object Methods	Description
<code>group(num=0)</code>	This method returns entire match (or specific subgroup num)
<code>groups()</code>	This method returns all matching subgroups in a tuple (empty if there weren't any)

Example:

```
#!/usr/bin/python
import re

line = "Cats are smarter than dogs"

matchObj = re.match( r'(.*) are (.*?) .*', line, re.M|re.I)

if matchObj:
    print "matchObj.group() : ", matchObj.group()
    print "matchObj.group(1) : ", matchObj.group(1)
    print "matchObj.group(2) : ", matchObj.group(2)
else:
    print "No match!!"
```

When the above code is executed, it produces following result:

```
matchObj.group() : Cats are smarter than dogs
matchObj.group(1) : Cats
matchObj.group(2) : smarter
```

The search Function

This function searches for first occurrence of RE *pattern* within *string* with optional *flags*.

Here is the syntax for this function:

```
re.search(pattern, string, flags=0)
```

Here is the description of the parameters:

Parameter	Description
pattern	This is the regular expression to be matched.
string	This is the string, which would be searched to match the pattern anywhere in the string.
flags	You can specify different flags using bitwise OR (). These are modifiers, which are listed in the table below.

The `re.search` function returns a **match** object on success, **None** on failure. We would use `group(num)` or `groups()` function of **match** object to get matched expression.

Match Object Methods	Description
<code>group(num=0)</code>	This method returns entire match (or specific subgroup num)
<code>groups()</code>	This method returns all matching subgroups in a tuple (empty if there weren't any)

Example:

```
#!/usr/bin/python
import re

line = "Cats are smarter than dogs";

matchObj = re.match( r'(.*) are (.*?) .*', line, re.M|re.I)

if matchObj:
    print "matchObj.group() : ", matchObj.group()
    print "matchObj.group(1) : ", matchObj.group(1)
    print "matchObj.group(2) : ", matchObj.group(2)
else:
    print "No match!!"
```

When the above code is executed, it produces following result:

```
matchObj.group() : Cats are smarter than dogs
matchObj.group(1) : Cats
matchObj.group(2) : smarter
```

Matching vs Searching:

Python offers two different primitive operations based on regular expressions: **match** checks for a match only at the beginning of the string, while **search** checks for a match anywhere in the string (this is what Perl does by default).

Example:

```
#!/usr/bin/python
import re

line = "Cats are smarter than dogs";

matchObj = re.match( r'dogs', line, re.M|re.I)
if matchObj:
    print "match --> matchObj.group() : ", matchObj.group()
else:
    print "No match!!"

matchObj = re.search( r'dogs', line, re.M|re.I)
if matchObj:
    print "search --> matchObj.group() : ", matchObj.group()
else:
    print "No match!!"
```

When the above code is executed, it produces the following result:

```
No match!!
search --> matchObj.group() : dogs
```

Search and Replace:

Some of the most important **re** methods that use regular expressions is **sub**.

Syntax:

```
re.sub(pattern, repl, string, max=0)
```

This method replaces all occurrences of the RE *pattern* in *string* with *repl*, substituting all occurrences unless *max* provided. This method would return modified string.

Example:

Following is the example:

```
#!/usr/bin/python
import re

phone = "2004-959-559 # This is Phone Number"

# Delete Python-style comments
num = re.sub(r'#. *$', "", phone)
print "Phone Num : ", num

# Remove anything other than digits
num = re.sub(r'\D', "", phone)
print "Phone Num : ", num
```

When the above code is executed, it produces the following result:

```
Phone Num : 2004-959-559
Phone Num : 2004959559
```

Regular-expression Modifiers - Option Flags

Regular expression literals may include an optional modifier to control various aspects of matching. The modifiers are specified as an optional flag. You can provide multiple modifiers using exclusive OR (`|`), as shown previously and may be represented by one of these:

Modifier	Description
<code>re.I</code>	Performs case-insensitive matching.
<code>re.L</code>	Interprets words according to the current locale. This interpretation affects the alphabetic group (<code>\w</code> and <code>\W</code>), as well as word boundary behavior (<code>\b</code> and <code>\B</code>).
<code>re.M</code>	Makes <code>\$</code> match the end of a line (not just the end of the string) and makes <code>^</code> match the start of any line (not just the start of the string).
<code>re.S</code>	Makes a period (dot) match any character, including a newline.
<code>re.U</code>	Interprets letters according to the Unicode character set. This flag affects the behavior of <code>\w</code> , <code>\W</code> , <code>\b</code> , <code>\B</code> .
<code>re.X</code>	Permits "cuter" regular expression syntax. It ignores whitespace (except inside a set <code>[]</code> or when escaped by a backslash) and treats unescaped <code>#</code> as a comment marker.

Regular-expression patterns:

Except for control characters, (`+ ? . * ^ $ () [] { } | \`), all characters match themselves. You can escape a control character by preceding it with a backslash.

Following table lists the regular expression syntax that is available in Python:

Pattern	Description
<code>^</code>	Matches beginning of line.
<code>\$</code>	Matches end of line.
<code>.</code>	Matches any single character except newline. Using <code>m</code> option allows it to match newline as well.
<code>[...]</code>	Matches any single character in brackets.
<code>[^...]</code>	Matches any single character not in brackets
<code>re*</code>	Matches 0 or more occurrences of preceding expression.
<code>re+</code>	Matches 1 or more occurrence of preceding expression.
<code>re?</code>	Matches 0 or 1 occurrence of preceding expression.
<code>re{n}</code>	Matches exactly <code>n</code> number of occurrences of preceding expression.
<code>re{n,}</code>	Matches <code>n</code> or more occurrences of preceding expression.
<code>re{n, m}</code>	Matches at least <code>n</code> and at most <code>m</code> occurrences of preceding expression.
<code>a b</code>	Matches either <code>a</code> or <code>b</code> .
<code>(re)</code>	Groups regular expressions and remembers matched text.
<code>(?imx)</code>	Temporarily toggles on <code>i</code> , <code>m</code> , or <code>x</code> options within a regular expression. If in parentheses, only that area is affected.

(?-imx)	Temporarily toggles off i, m, or x options within a regular expression. If in parentheses, only that area is affected.
(?: re)	Groups regular expressions without remembering matched text.
(?imx: re)	Temporarily toggles on i, m, or x options within parentheses.
(?-imx: re)	Temporarily toggles off i, m, or x options within parentheses.
(?#...)	Comment.
(?= re)	Specifies position using a pattern. Doesn't have a range.
(?! re)	Specifies position using pattern negation. Doesn't have a range.
(?> re)	Matches independent pattern without backtracking.
\w	Matches word characters.
\W	Matches nonword characters.
\s	Matches whitespace. Equivalent to [\t\n\r\f].
\S	Matches nonwhitespace.
\d	Matches digits. Equivalent to [0-9].
\D	Matches nondigits.
\A	Matches beginning of string.
\Z	Matches end of string. If a newline exists, it matches just before newline.
\z	Matches end of string.
\G	Matches point where last match finished.
\b	Matches word boundaries when outside brackets. Matches backspace (0x08) when inside brackets.
\B	Matches nonword boundaries.
\n, \t, etc.	Matches newlines, carriage returns, tabs, etc.
\1...\9	Matches nth grouped subexpression.
\10	Matches nth grouped subexpression if it matched already. Otherwise refers to the octal representation of a character code.

REGULAR-EXPRESSION EXAMPLES

Literal characters:

Example	Description
python	Match "python".

Character classes:

Example	Description
[Pp]ython	Match "Python" or "python"
rub[ye]	Match "ruby" or "rube"
[aeiou]	Match any one lowercase vowel
[0-9]	Match any digit; same as [0123456789]
[a-z]	Match any lowercase ASCII letter
[A-Z]	Match any uppercase ASCII letter
[a-zA-Z0-9]	Match any of the above
[^aeiou]	Match anything other than a lowercase vowel
[^0-9]	Match anything other than a digit

Special Character Classes:

Example	Description
.	Match any character except newline
\d	Match a digit: [0-9]
\D	Match a nondigit: [^0-9]
\s	Match a whitespace character: [\t\r\n\f]
\S	Match nonwhitespace: [^\t\r\n\f]
\w	Match a single word character: [A-Za-z0-9_]
\W	Match a nonword character: [^A-Za-z0-9_]

Repetition Cases:

Example	Description
ruby?	Match "rub" or "ruby": the y is optional
ruby*	Match "rub" plus 0 or more ys
ruby+	Match "rub" plus 1 or more ys
\d{3}	Match exactly 3 digits
\d{3,}	Match 3 or more digits
\d{3,5}	Match 3, 4, or 5 digits

Nongreedy repetition:

This matches the smallest number of repetitions:

Example	Description
<code><.*></code>	Greedy repetition: matches "<python>perl"
<code><.*?></code>	Nongreedy: matches "<python>" in "<python>perl"

Grouping with parentheses:

Example	Description
<code>\D\d+</code>	No group: + repeats \d
<code>(\D\d)+</code>	Grouped: + repeats \D\d pair
<code>([Pp]ython(,)?)+</code>	Match "Python", "Python, python, python", etc.

Backreferences:

This matches a previously matched group again:

Example	Description
<code>([Pp])ython&\1ails</code>	Match python&pails or Python&Pails
<code>(["'])[^\1]*\1</code>	Single or double-quoted string. \1 matches whatever the 1st group matched. \2 matches whatever the 2nd group matched, etc.

Alternatives:

Example	Description
<code>python perl</code>	Match "python" or "perl"
<code>rub(y le)</code>	Match "ruby" or "ruble"
<code>Python(!+ \?)</code>	"Python" followed by one or more ! or one ?

Anchors:

This needs to specify match position.

Example	Description
<code>^Python</code>	Match "Python" at the start of a string or internal line
<code>Python\$</code>	Match "Python" at the end of a string or line
<code>\APython</code>	Match "Python" at the start of a string
<code>Python\Z</code>	Match "Python" at the end of a string
<code>\bPython\b</code>	Match "Python" at a word boundary

<code>\brub\b</code>	<code>\B</code> is nonword boundary: match "rub" in "rube" and "ruby" but not alone
<code>Python(?!)</code>	Match "Python", if followed by an exclamation point
<code>Python(?!)</code>	Match "Python", if not followed by an exclamation point

Special syntax with parentheses:

Example	Description
<code>R(?#comment)</code>	Matches "R". All the rest is a comment
<code>R(?i)uby</code>	Case-insensitive while matching "uby"
<code>R(?i:uby)</code>	Same as above
<code>rub(?:y le)</code>	Group only without creating <code>\1</code> backreference